

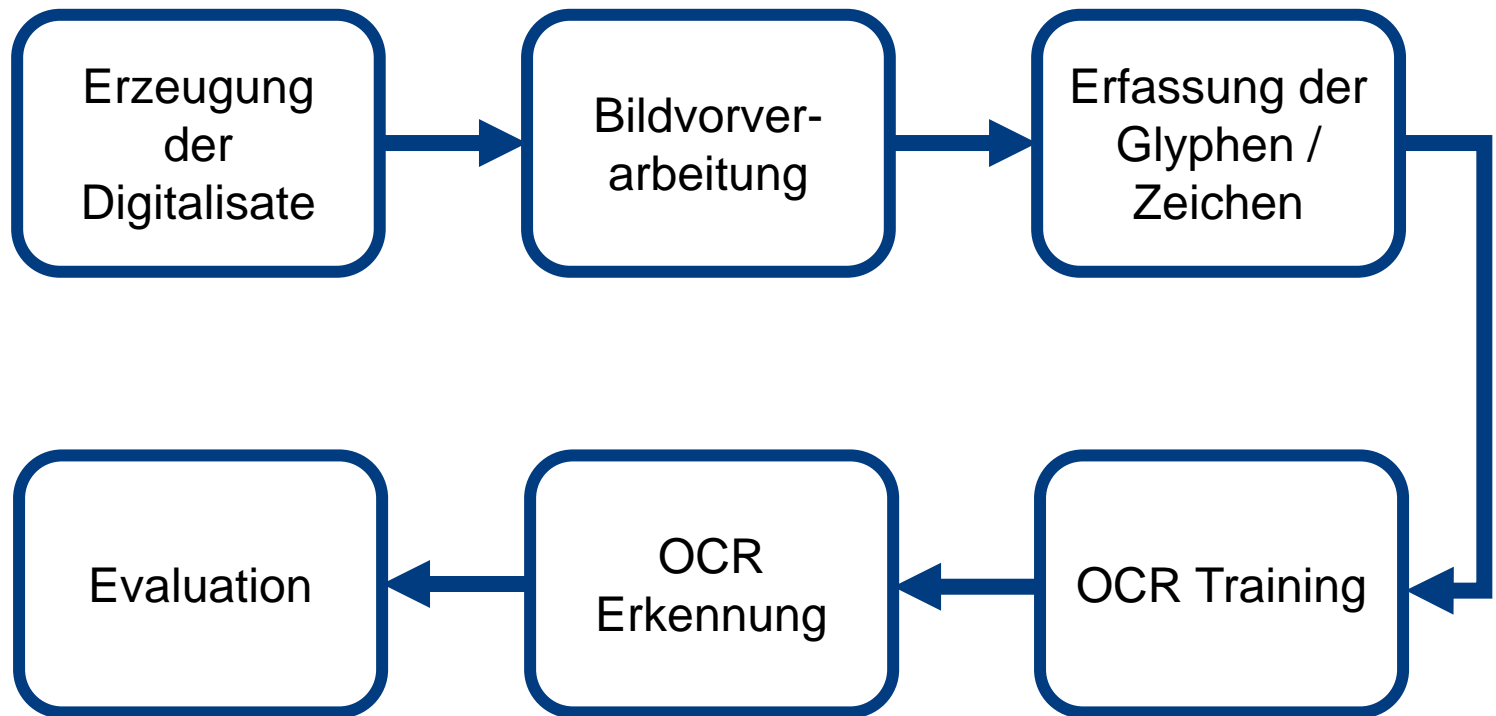
**<philtag n="13"/>**

# **OCR Workshop**

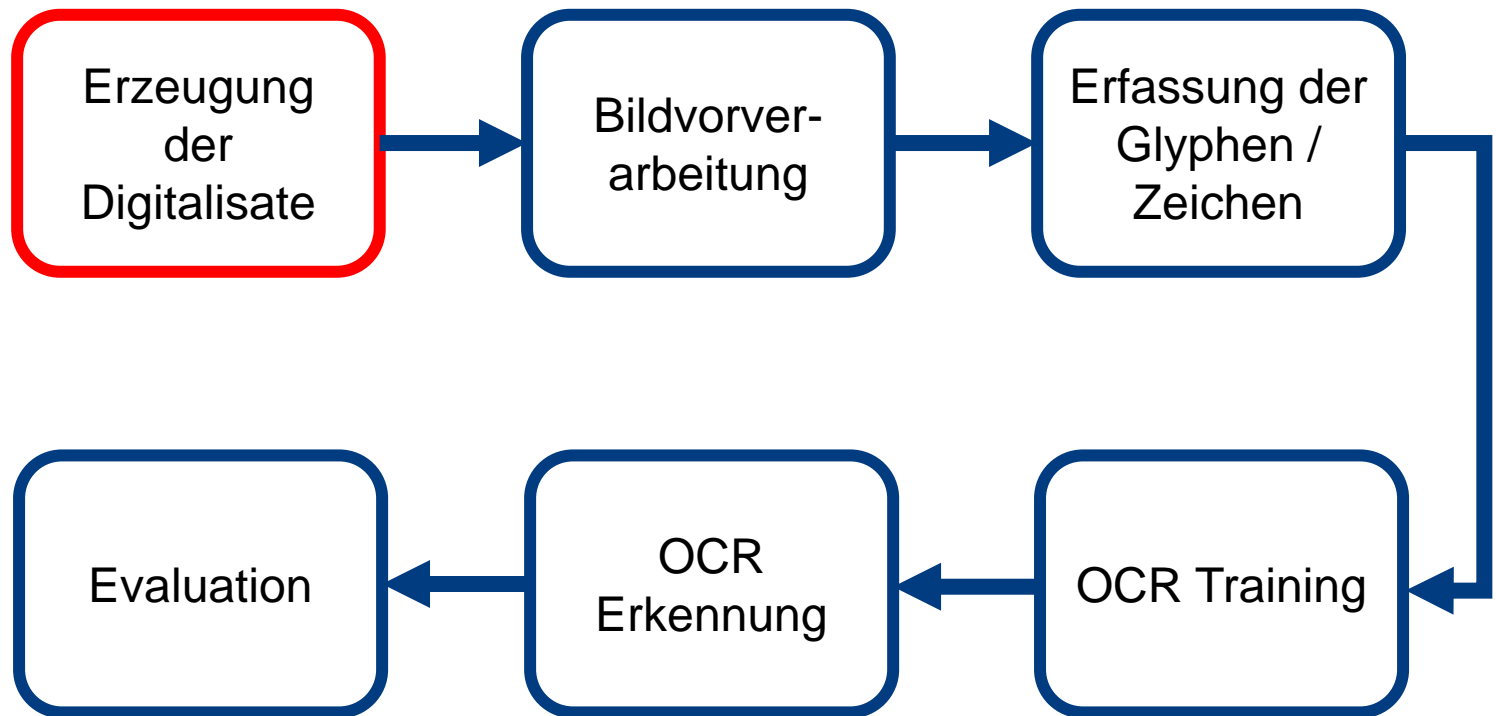
# Tagesordnung

- 13:40 Digitalisate
- 14:00 Glyphen
- 14:15 Glyph Miner
- 14:45 Pause
- 15:00 Aletheia
- 15:45 Franken+
- 16:30 VietOCR

# OCR Workflow



# OCR Workflow

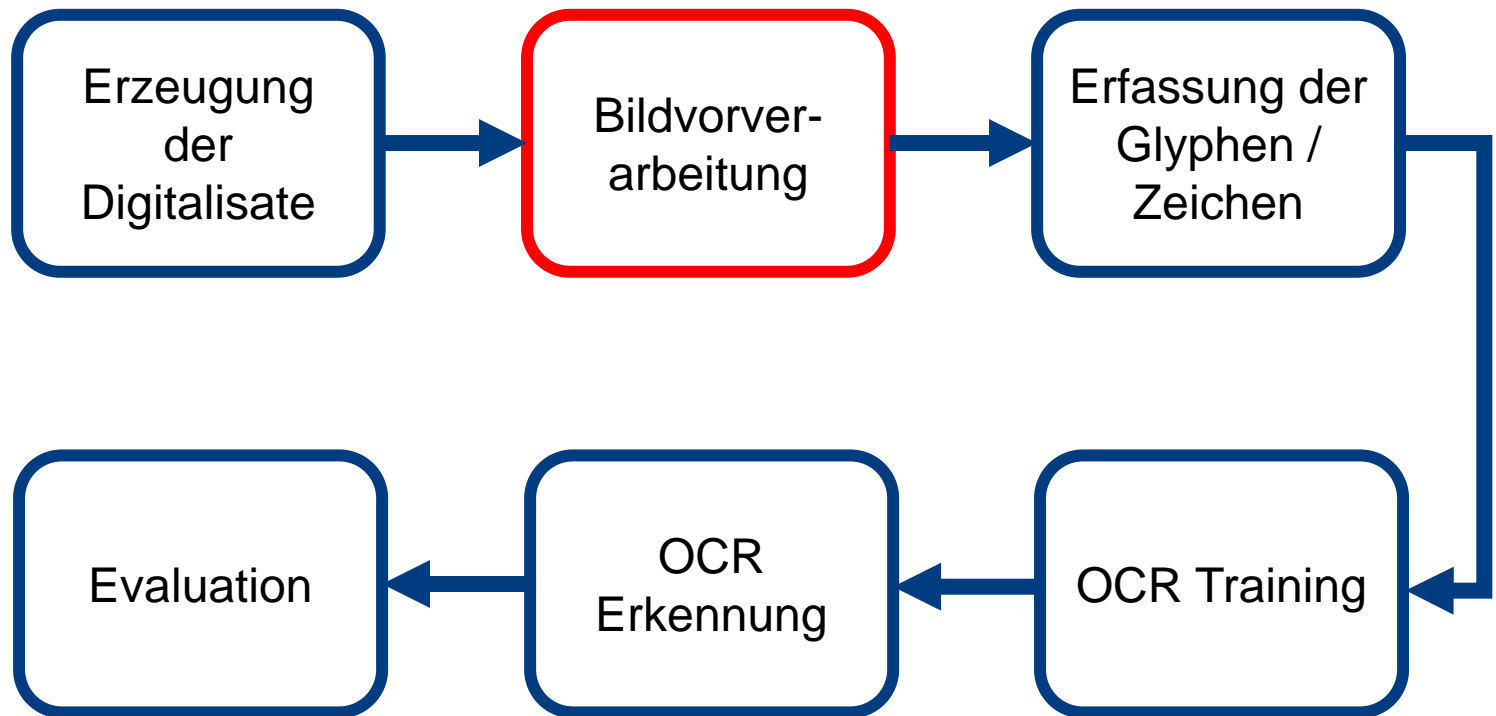


# Erstellung der Digitalisate

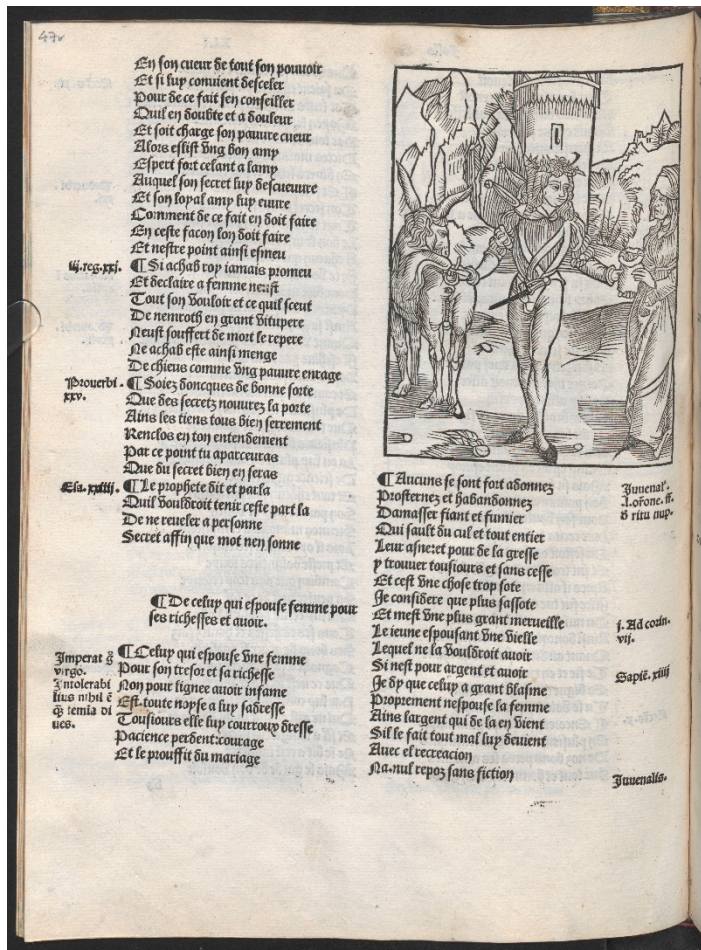
## Anforderungen:

- Scanauflösung: mind. 300 DPI
- Ausgerichtete Aufnahme
- Verzerrungsfreie Aufnahme:
  - Kein Umbug im Buchfalz
  - Keine Wellen im Papier

# OCR Workflow



# Bildvorverarbeitung



## Problemstellung:

- Automatische Entzerrung und Rotation
- Automatische Binarisierung
- Entfernung von Bildstörungen
- Segmentierung komplexer Layouts
- Beibehaltung der Lesereihenfolge

# Bildvorverarbeitung: Manuell

Bildbearbeitungsprogramme,

z.B. [IrfanView](#)

- Vorteile: kostenlos, einfach zu bedienen, volle Kontrolle durch den Nutzer
- Nachteile: zeitaufwändig, keine Entzerrung und Rotation



# Bildvorverarbeitung: Automatisch #1

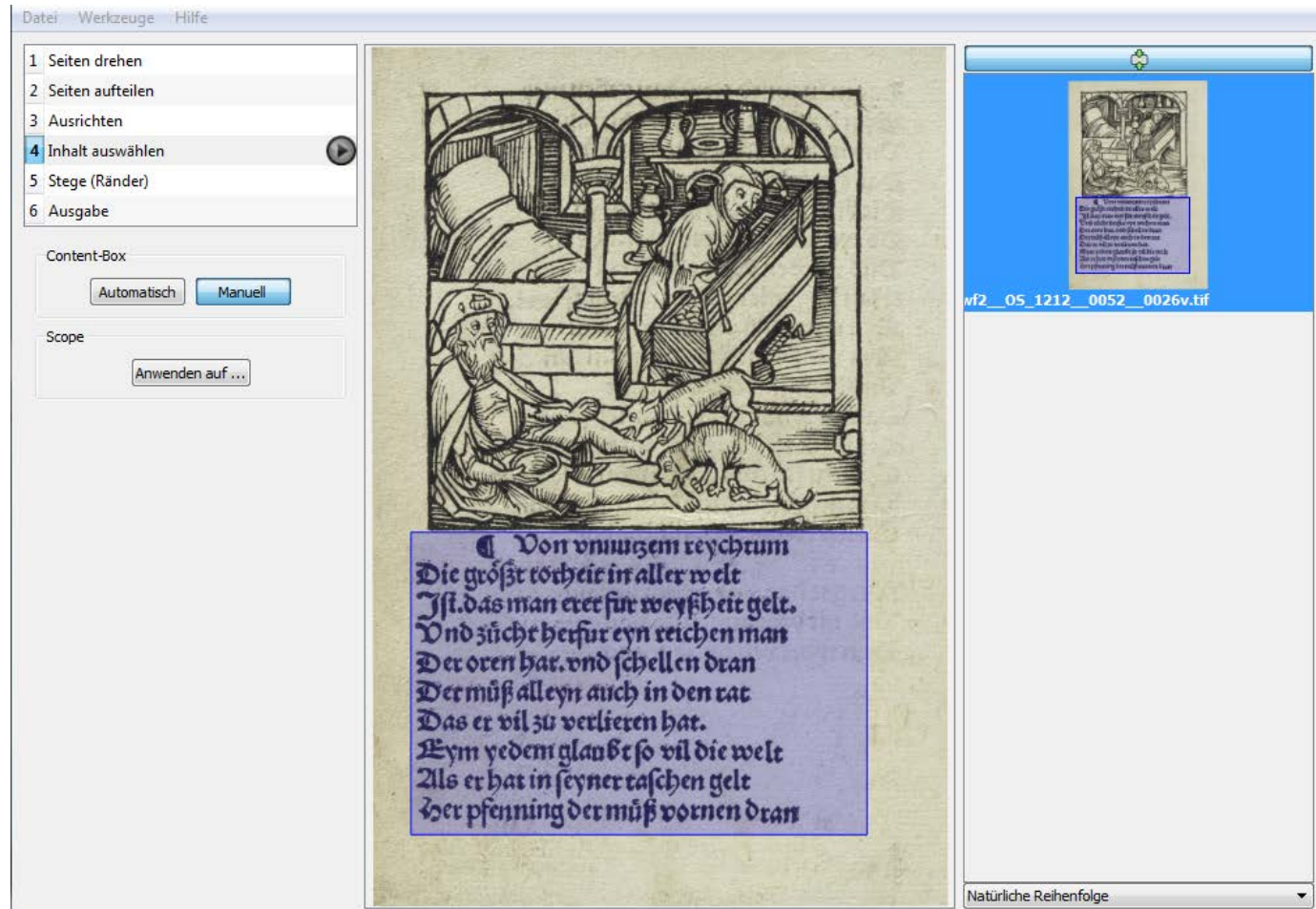
Nicht-kommerziell, z.B. [ScanTailor](#)

- Vorteile:
  - OpenSource und Plattform-unabhängig
  - halbautomatische Entzerrung, Rotation und Beschnitt
  - Binarisierung mit Entfernung von Bildstörungen wie Flecken und Rauschen

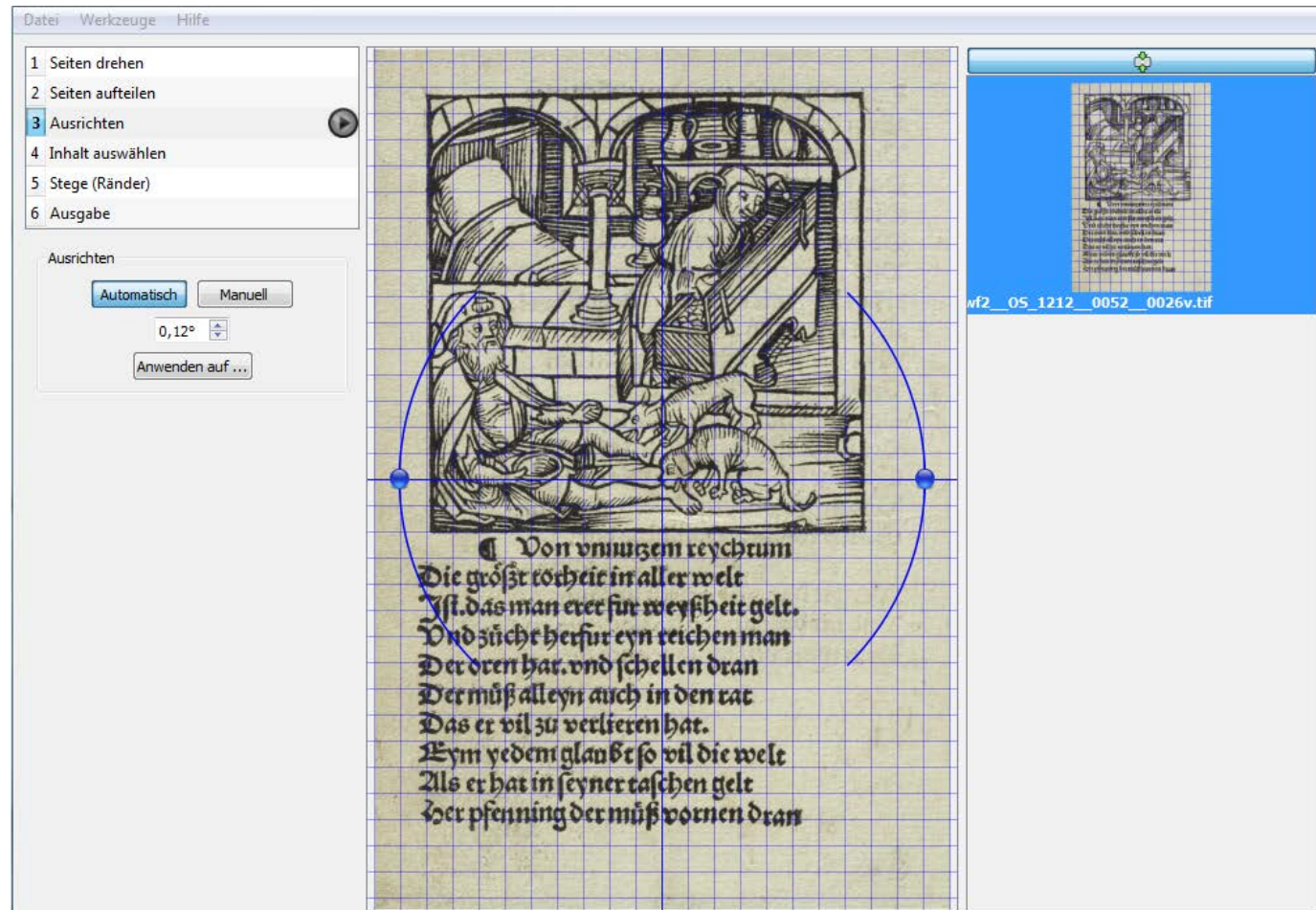
# Bildvorverarbeitung: Automatisch #2

- Nachteile:
  - Automatik-Funktionen liefern teilweise falsche Ergebnisse
  - Komplexere Layouts schwer erfassbar
  - Nachkontrolle durch den Nutzer notwendig

# Bildvorverarbeitung: ScanTailor #1

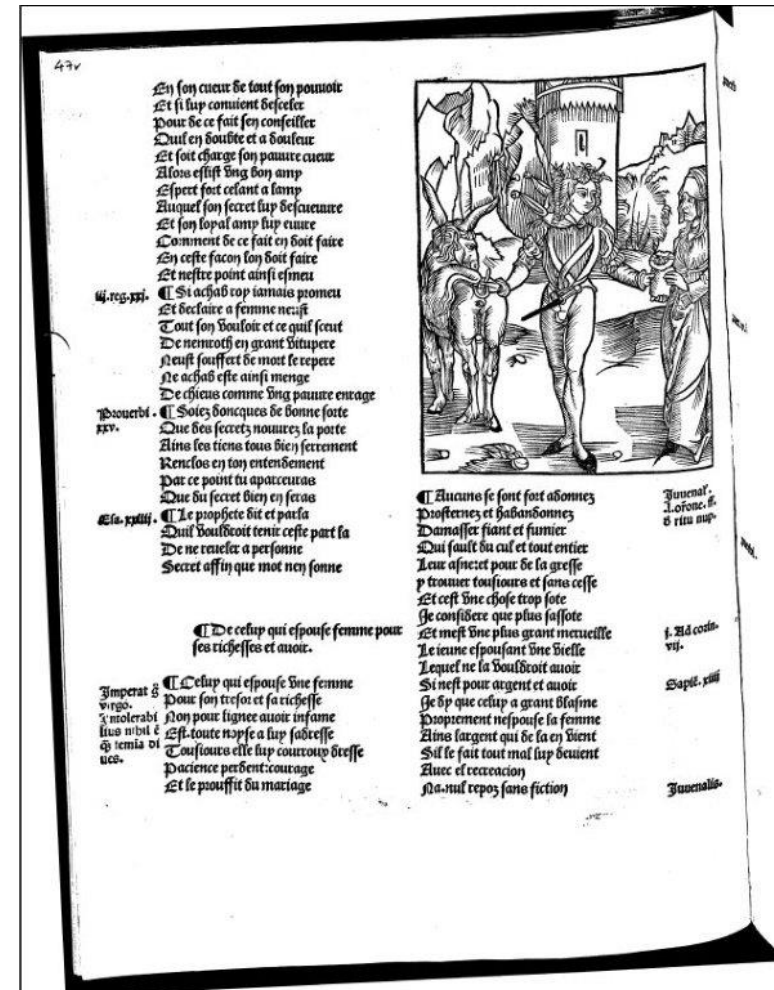
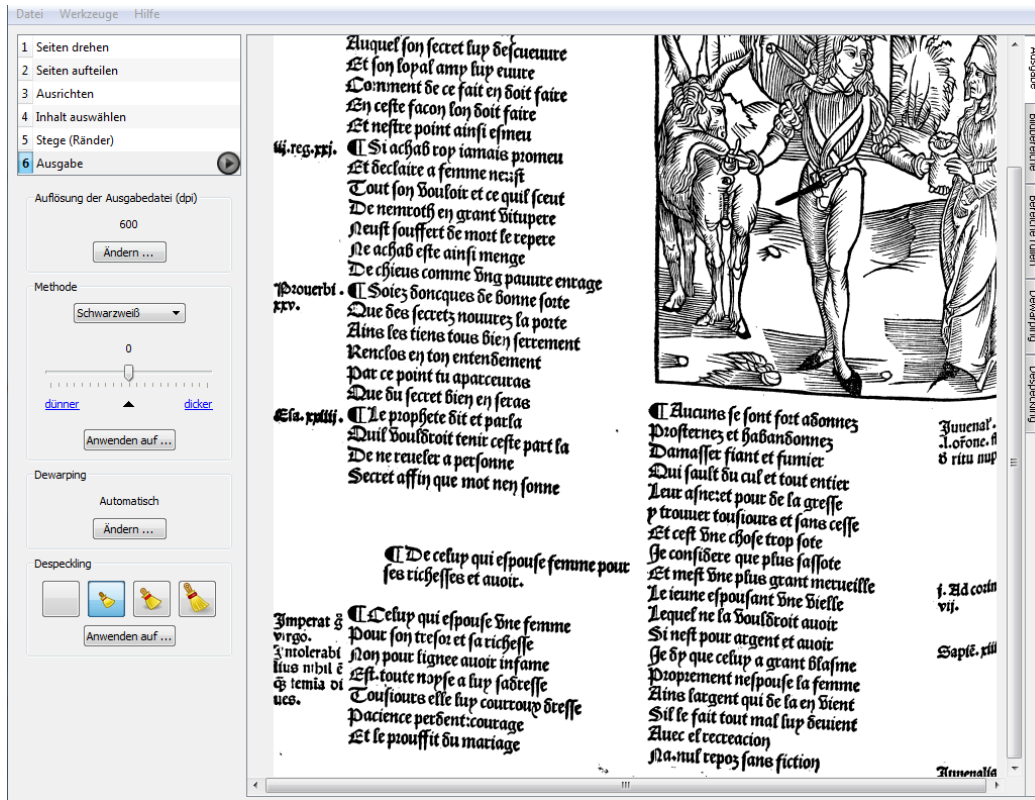


# Bildvorverarbeitung: ScanTailor #2





# Bildvorverarbeitung: ScanTailor vs. Abbyy FineReader

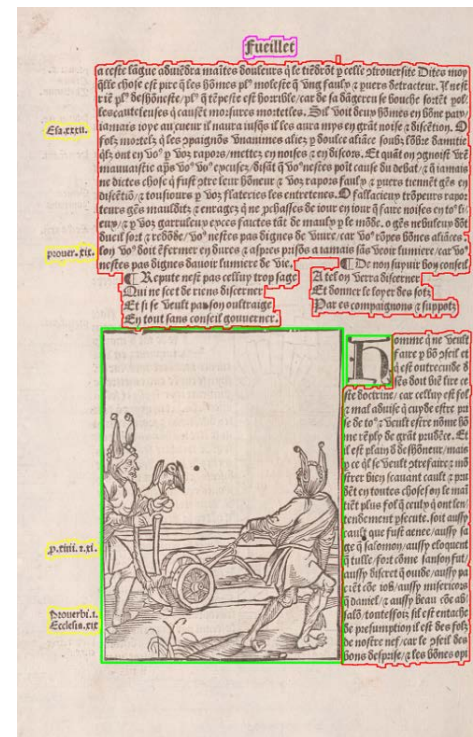
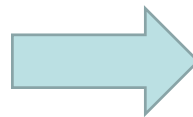
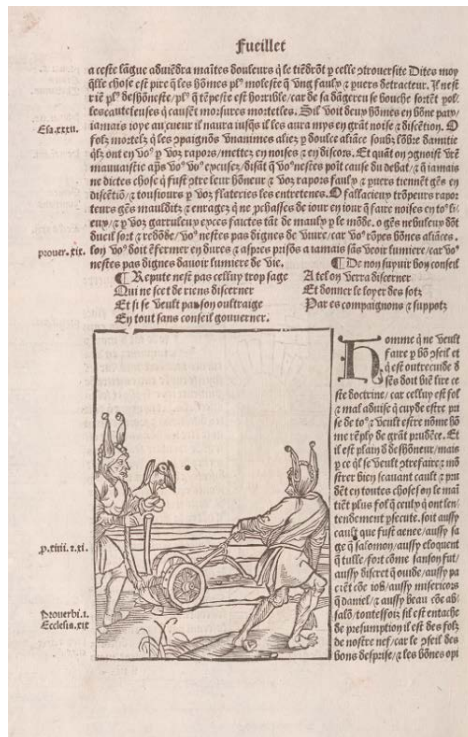


# Segmentierung von Drucken

Christian Reul

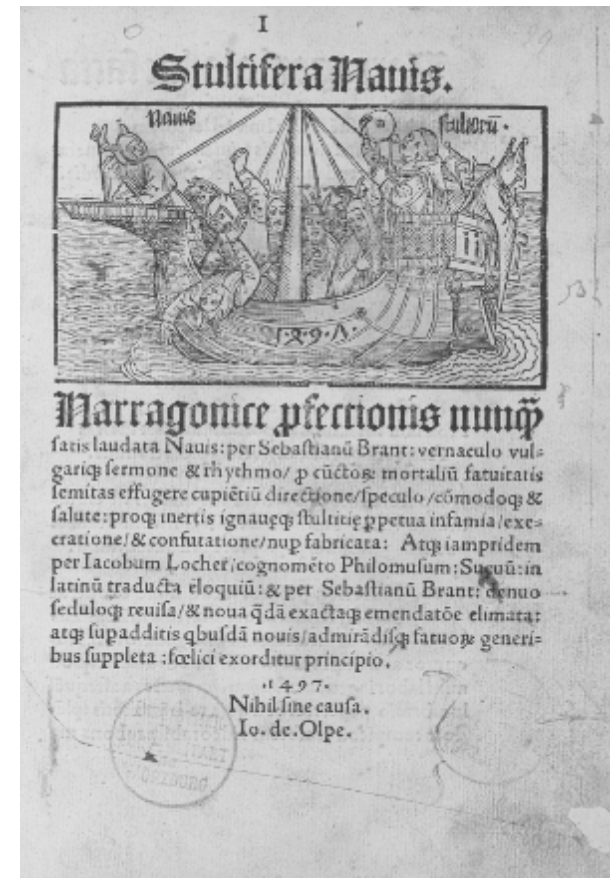
Prof. Dr. Frank Puppe

Lehrstuhl für Künstliche Intelligenz und Angewandte Informatik

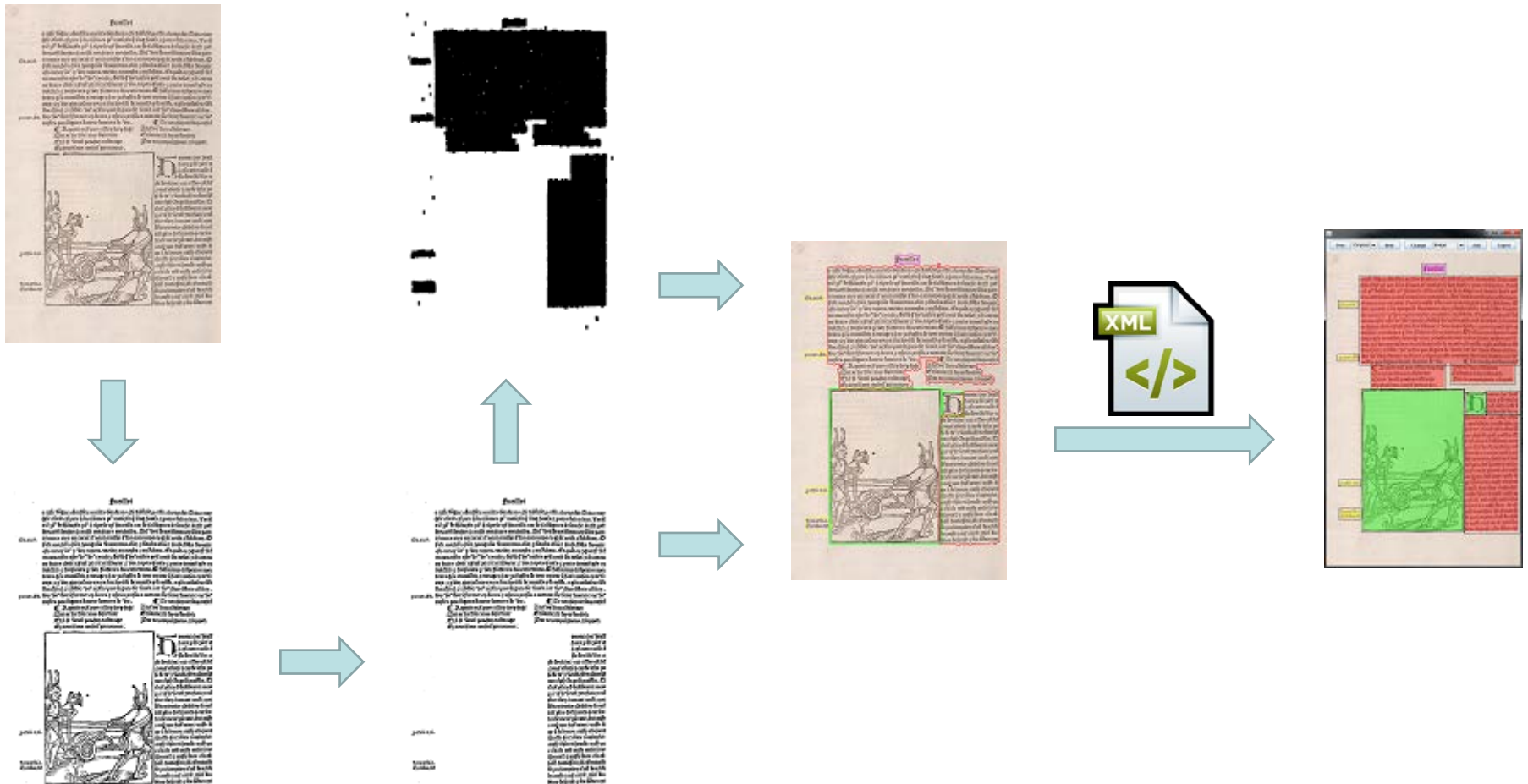


# Motivation und Konzept

- Einfache, schnelle Segmentierung ausgewählter Werke.
  - Verschiedene Narrenschiff-Ausgaben.
  - Der neue/praktische Schulmann.
- Anwendung grundlegender Methoden aus der Bildverarbeitung.
- Aufzeigen und Lösen/Umgehen der Probleme, die bei einem simplen Ansatz auftreten.
- Tool zur einfachen Korrektur einer fehlerhaften Lösung.



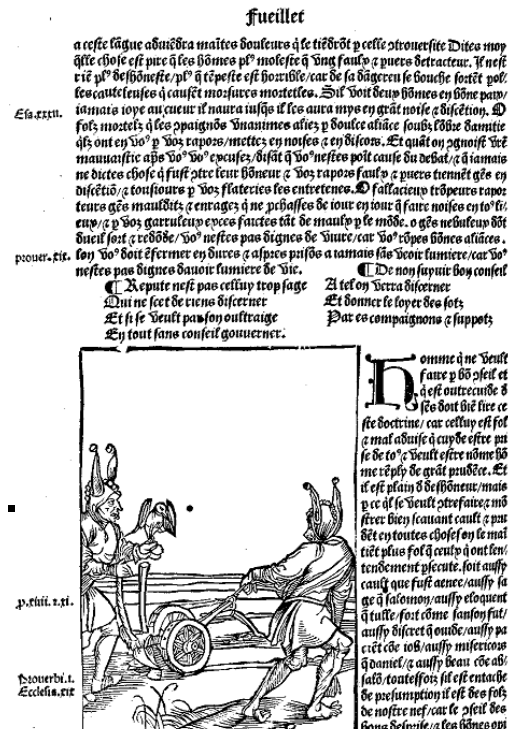
# Schematischer Ablauf





# Vorverarbeitung und Bilddetektion

- Zuschneiden und Ausrichten (Uni-Bibliothek).
- (Temporäre) Skalierung auf einheitliche Größe.
- Binärisierung nach Otsu.
- Bilddetektion und -entfernung durch Konturfindung.



Binärdarstellung.

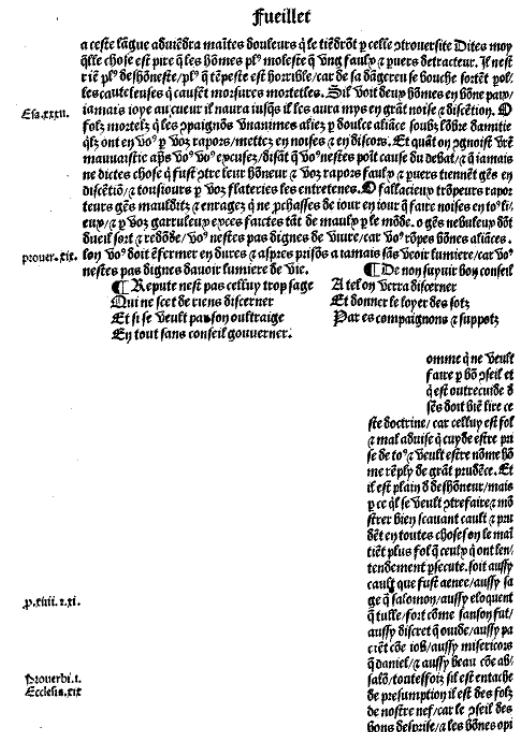


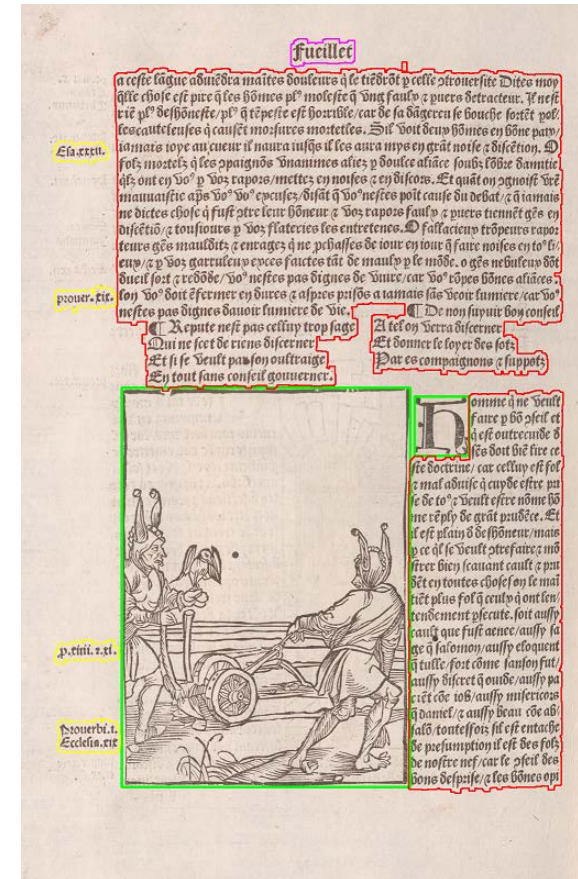
Bild und Initialie entfernt.

# Klassifikation von Textabschnitten

- Verbinden von Regionen durch Schließen von Zwischenräumen.
- Klassifikation anhand von Größe und Position.
- Kategorien:
  - Paragraph
  - Überschrift
  - Marginalie
  - Seitennummer
  - ...



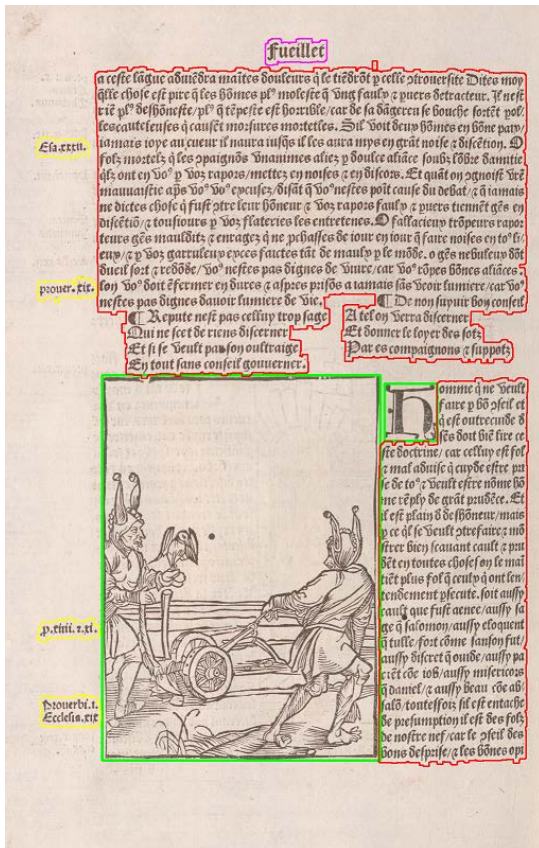
Verbundene Textregionen.



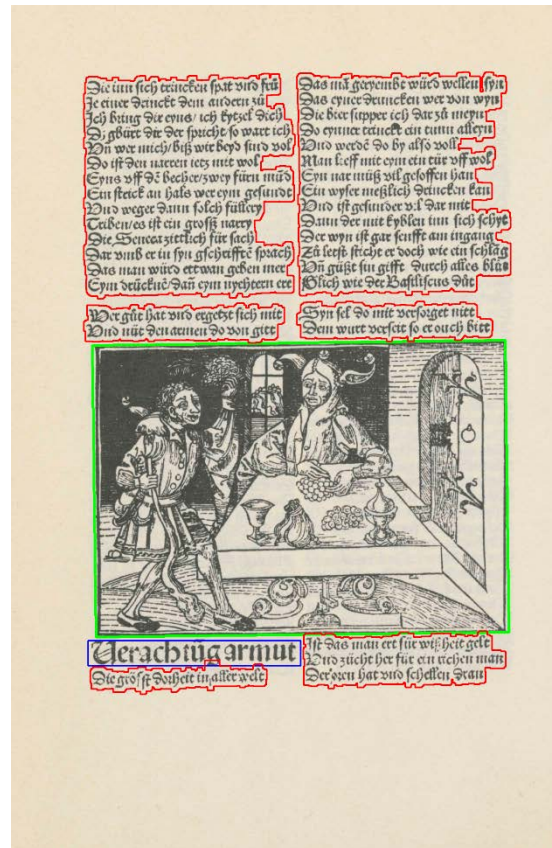
Finale Segmentierung.



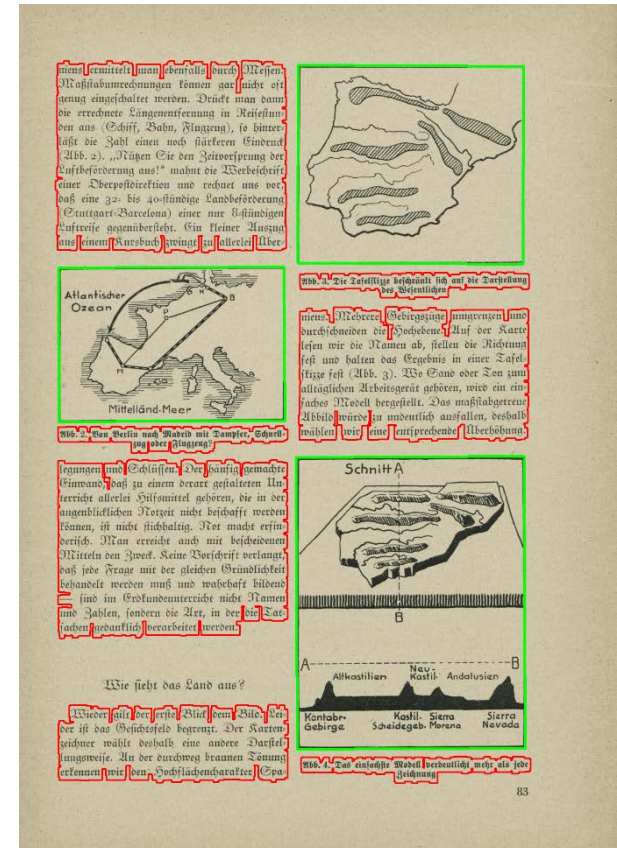
# Beispielesegmentierungen



Narrenschiff GW5060, Lyon 1499.



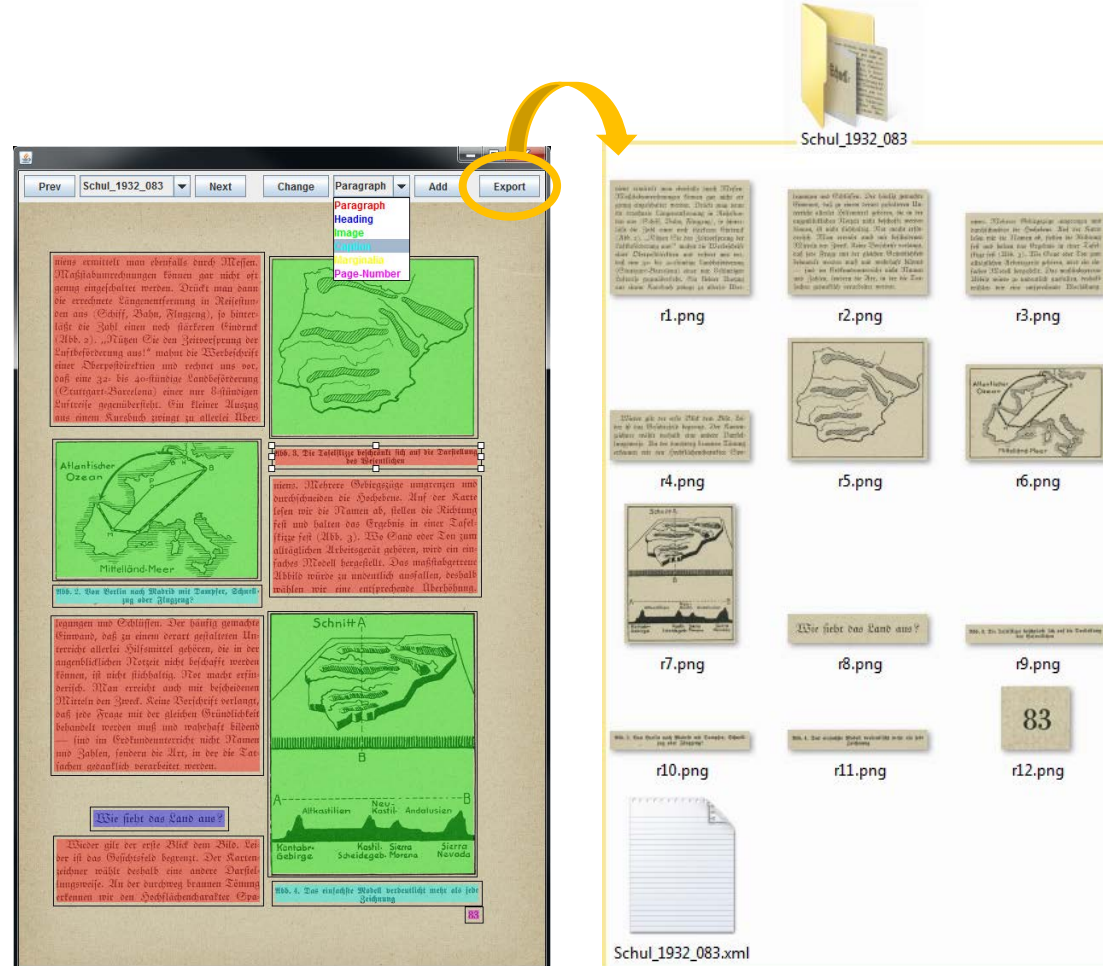
Narrenschiff GW5048, Straßburg 1494.



Der praktische Schulmann, 1932 - Heft 6.

# Tool zur manuellen Nachbearbeitung

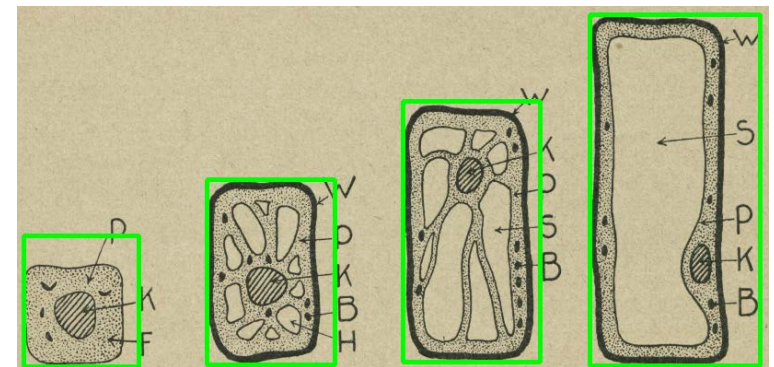
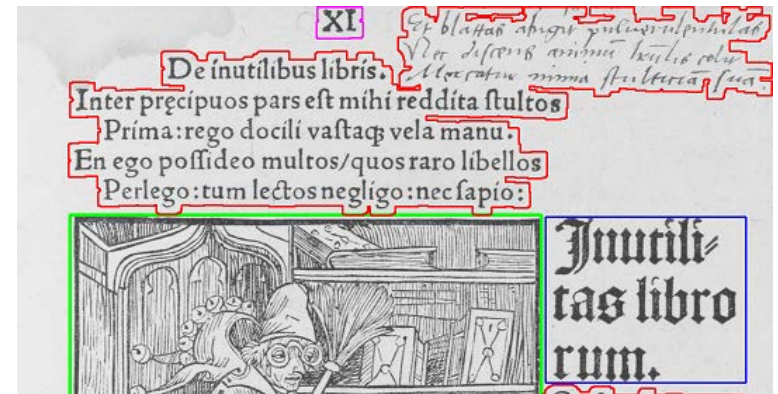
- Anzeigen des Bildes und der zugehörigen Segmentierung (PageXML).
- Löschen, Hinzufügen und Ändern von Regionen.
- Abspeichern der einzelnen Regionen und der zugehörigen PageXML-Datei.





# Probleme und Lösungsansätze bei der Segmentierung von Drucken I

- Verunreinigungen der Ausgangsbilder.
  - Allgemeingültige Lösung schwierig.
  - Oft durch simple Nachbearbeitung behebbar.
- Mit diesem Ansatz keine verlässliche Bilddetektion möglich.
  - Bisher lediglich Orientierung an der äußeren Begrenzung des Bildes.
  - Fehldetektionen können die Segmentierung der gesamten Seite quasi unbrauchbar machen.



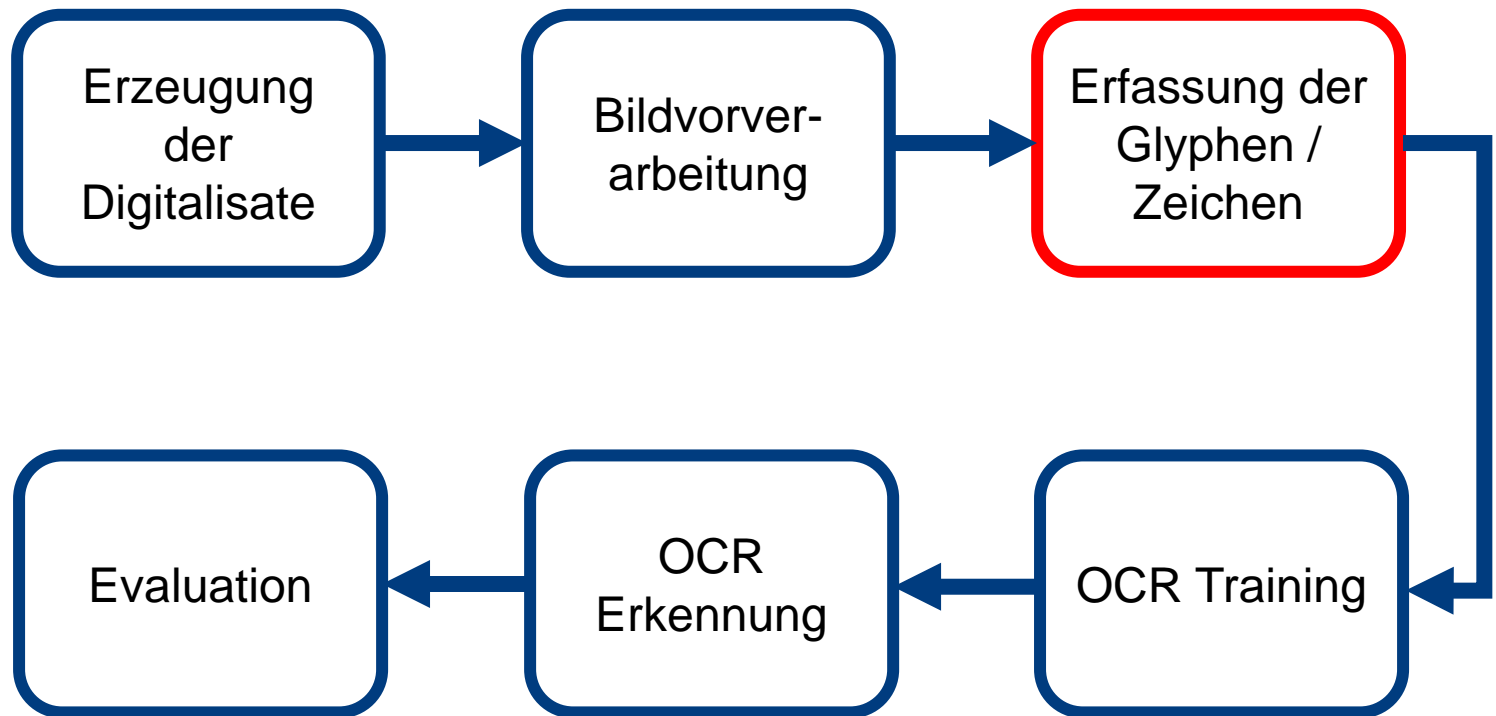
# Probleme und Lösungsansätze bei der Segmentierung von Drucken II

- Beschränkte Generalisierbarkeit durch starke Parametrisierung.
  - Nachteil: Manueller Aufwand.
    - Unterstützung denkbar durch:
      - (Halb)automatische Parameteroptimierung durch iterativen Ansatz.
      - Interaktive Festlegung (z. B. Konfiguration anhand von Beispielseite).
  - Vorteil: Individuell anpassbar.

## Fazit:

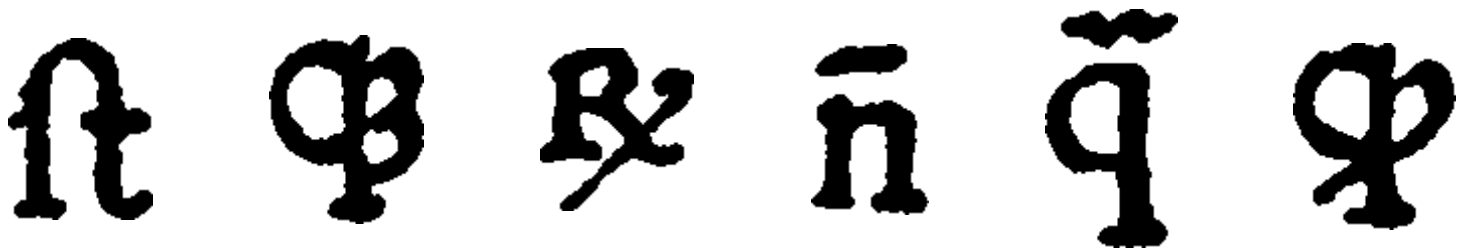
- Sehr fordernde Aufgabe, obwohl für den Menschen trivial.
- Selbst ein simpler Ansatz kann gute Ergebnisse liefern.
- Qualität stark abhängig vom zu segmentierenden Werk.
- Universallösung mit simplen Methoden unmöglich.

# OCR Workflow



# Glyphen #1

- Frühe Drucke beinhalten häufig eine Vielzahl an Sonderzeichen:
  - Ligaturen
  - spezielle Abkürzungen
  - Vokale mit Diakritika





## Glyphen #2

- Die Sonderzeichen sollten für eine 1:1 Übertragung bei der OCR weitestgehend erhalten bleiben:
  - Kein Vorgreifen auf eine wissenschaftliche Auswertung
  - Abkürzungen können kontextabhängig sein
  - Ligaturen können Wortgrenzen verdeutlichen

# Glyphen - MUFI

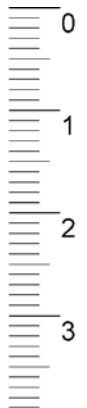
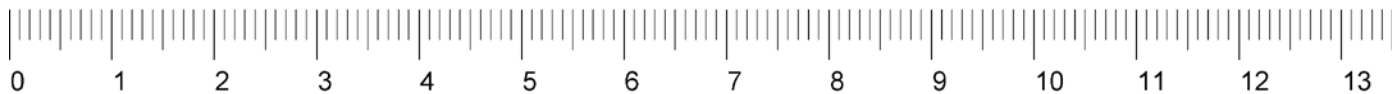
- Medieval Unicode Font Initiative (MUFI):
  - Internationales Projekt zur Kodierung spezieller lateinischer Zeichen im Unicode Standard
  - Verwendet die sog. Private Use Area des Unicode-Standards

# Glyphen - Typenrepertorium

- Typenrepertorium der Wiegendrucke:  
– beinhaltet die Typeninventare von ca. 2000 Offizinen (nicht vollständig)

A B C D E F G H I K L M N O P Q Q<sub>z</sub> R S T V W X  
a ā b c d e ē f f f l g h i k l l m n ñ o ō œ p p̄ p p q q̄ q q̄  
r r̄ s s f f t t t u ū v x y z 9 & } ! ( ) / = : • 1 4 7 9

Type 1.



# Glyphen - Typeninventar

- Ziel: Erfassung (möglichst) aller Schriftarten und Zeichen einer Offizin bzw. eines Druckers
- Beispiel: Bergmann von Olpe (Basel)
  - Type 1:109R, 2:180G, 3:77R, 4:220G, 5:109G

# Glyphen – Übertragbarkeit #1

- Annahme: Wurde die OCR für alle Typen eines Druckers trainiert, können alle Werke dieser Offizin per OCR verarbeitet werden.
- Beispiel:
  - 5 von 8 Typen der Offizin Bergmann von Olpe können mit nur 2 Werken trainiert werden.

# Glyphen – Übertragbarkeit #2

Enerat exactū post plurima sæcula tēpus:  
 Quo deus hūana vīsurus ymagīe terras.  
 Eligeret sacram tota de gente puellam:  
 Quę sibi mirificum casto de sanguīne corpus.  
 Texeret. & puero matris pia munera ferret.  
 Qua fatus humanis lumen pręferre tenebris  
 Vellet. & erratis varijs imponere finem.  
 Hanc sua virginitas .niuei sanctiq; pudoris  
 Inuiolatus amor .gratam fecere tonanti.  
 Letus ab ęthereis ad eam mox nuncius oris  
 Mittitur .vt summi ferret mandata parentis.  
 Pone metus (inquit) superis gratissima virgo.  
 Conceptura nouum nullo de femine fœtum:  
 Et prolem paritura dei: materq; futura es  
 Stīpīs olīmpiacę: tactus non passa viriles.  
 Et deus ipse tuam veniens labetur in alium:

Enerat~~o~~exactū po~~ft~~ p~~I~~urima læcula tēpus :  
 Quo deus hūan~~av~~īsurus ymagīe terras .  
 Eligeret sacram tota de gente puellam:  
 Quę ~~fl~~bi miri~~fl~~cum ca~~ft~~o de sanguine corpus .  
 Texeret. & .puero matris pia munera ferret.  
 Qua fatus humanis lumen pręferre tenebris  
 Vel~~g~~et. & erratis varil~~is~~ imponere finem .  
 Hanc sua virginitas .niuei sanctiq; pudoris  
 Inuiolatus amor .gratam fecere tonanti.  
 Letus ab ęthereis ad eam mox nuncius oris  
 Mittitur .vt summi ferret mandata parentis .  
~~p~~one metus (inquit)superis gratissima virgo .  
 Conceptura nouum nullo de femine fœtum:  
 Et prolem paritura dei : materq; futura es  
 Stū~~r~~.pīs olīmpiacę: tactus non pa~~ff~~a viriles.  
 Et deus ipse tuam veniens labetur in alium:

# Glyph Miner

- Ziel: Effiziente Erfassung von Glypheninventaren (mit Beispielen)

Character "g"



- Entwickelt an der Uni Würzburg  
<https://github.com/benedikt-budig/glyph-miner>

# Aletheia

- Entwickelt von PRImA (*Pattern Recognition & Image Analysis Research Lab*)
- University of Salford, Manchester
- <http://www.primaresearch.org/tools/Aletheia>



- Ziel: PAGE-XML (2010)

```

1 <Page imageFilename="GW5042__DE-20_FLDE-Swf2__OS_1212__0024__0011v_bw.tif" imageWidth="2038" imageHeight="3025">
2   <TextRegion id="r1" type="paragraph">
3     <Coords>
9     <TextLine id="l2">
10    <Coords>
16    <Word id="w3">
17      <Coords>
23      <Glyph id="c4">
24        <Coords>
25          <Point x="315" y="200"/>
26          <Point x="315" y="285"/>
27          <Point x="417" y="285"/>
28          <Point x="416" y="200"/>
29        </Coords>
30        <TextEquiv>
31          <PlainText/>
32          <Unicode>G</Unicode>
33        </TextEquiv>
34      </Glyph>
35      <Glyph id="c5">
36        <Coords>
37          <Point x="420" y="219"/>
38          <Point x="420" y="275"/>
39          <Point x="456" y="275"/>
40          <Point x="456" y="219"/>

```

# Aletheia

## Lite


- Binarisierung
- Schmutzentfernung
- Manuelle Zuweisung

## Pro

- Binarisierung
- Schmutzentfernung
- Automatische Erkennung (Regionen, Zeilen, Wörtern, Glyphen)

### Lite

**Aletheia Lite** offers essential features for viewing and editing page layout files. It supports various file formats including PAGE, ALTO and ABBYY FineReader XML. It comes with a **free and unlimited** licence for registered users (personal research).

 Login/Register to get licence key

### Pro

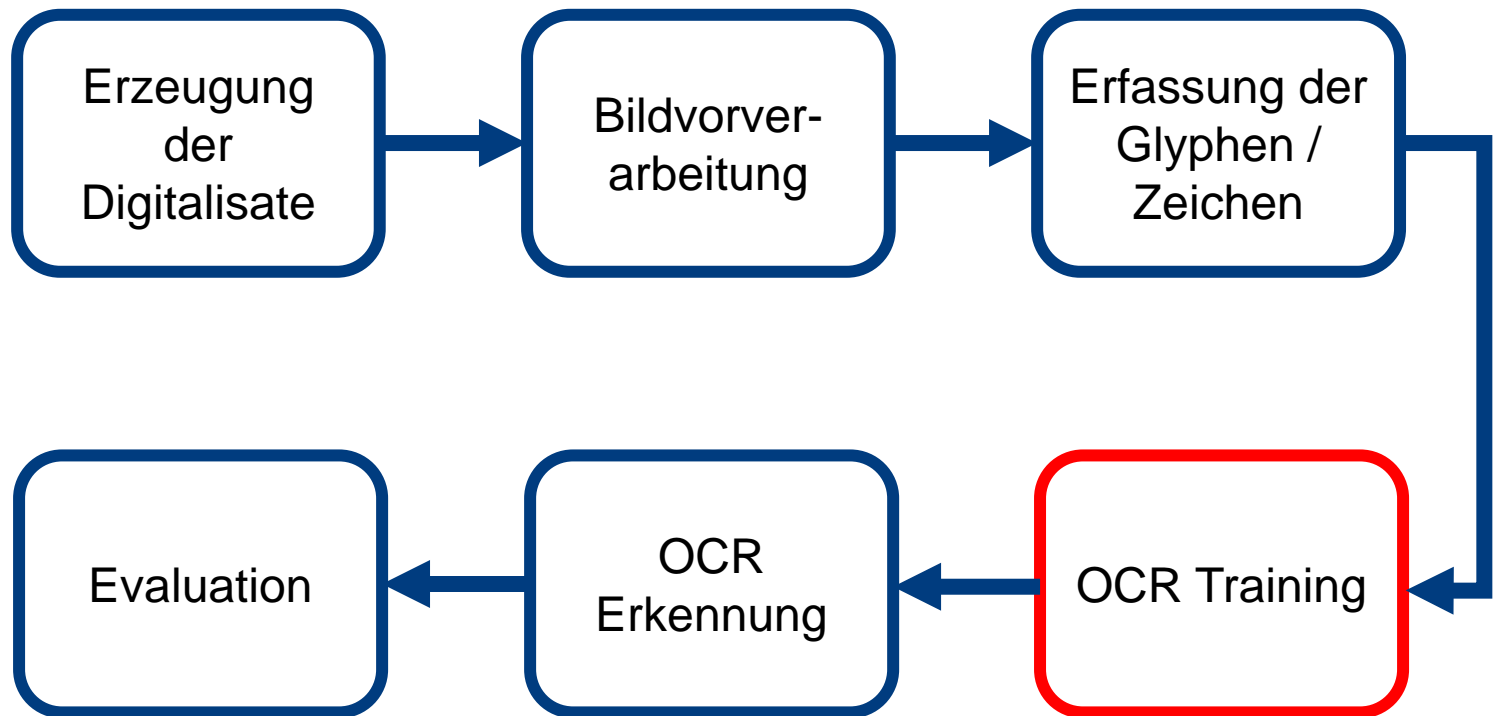
**Aletheia Pro** is a complete document analysis system offering a wide range of fully automated and assistive features including page recognition and OCR (Tesseract). The **complete product** is available to try for free for a limited period.

Free trial

Get Aletheia Pro

<http://www.primaresearch.org/tools/Aletheia/Editions>

# OCR Workflow



# Franken+

- Entwickelt vom eMOP *(Early Modern OCR Project)*
- Texas A&M University
- <http://emop.tamu.edu/outcomes/Franken-Plus>

## Voraussetzungen:

- .NET Framework
- MySQL
- Tesseract

## Material:

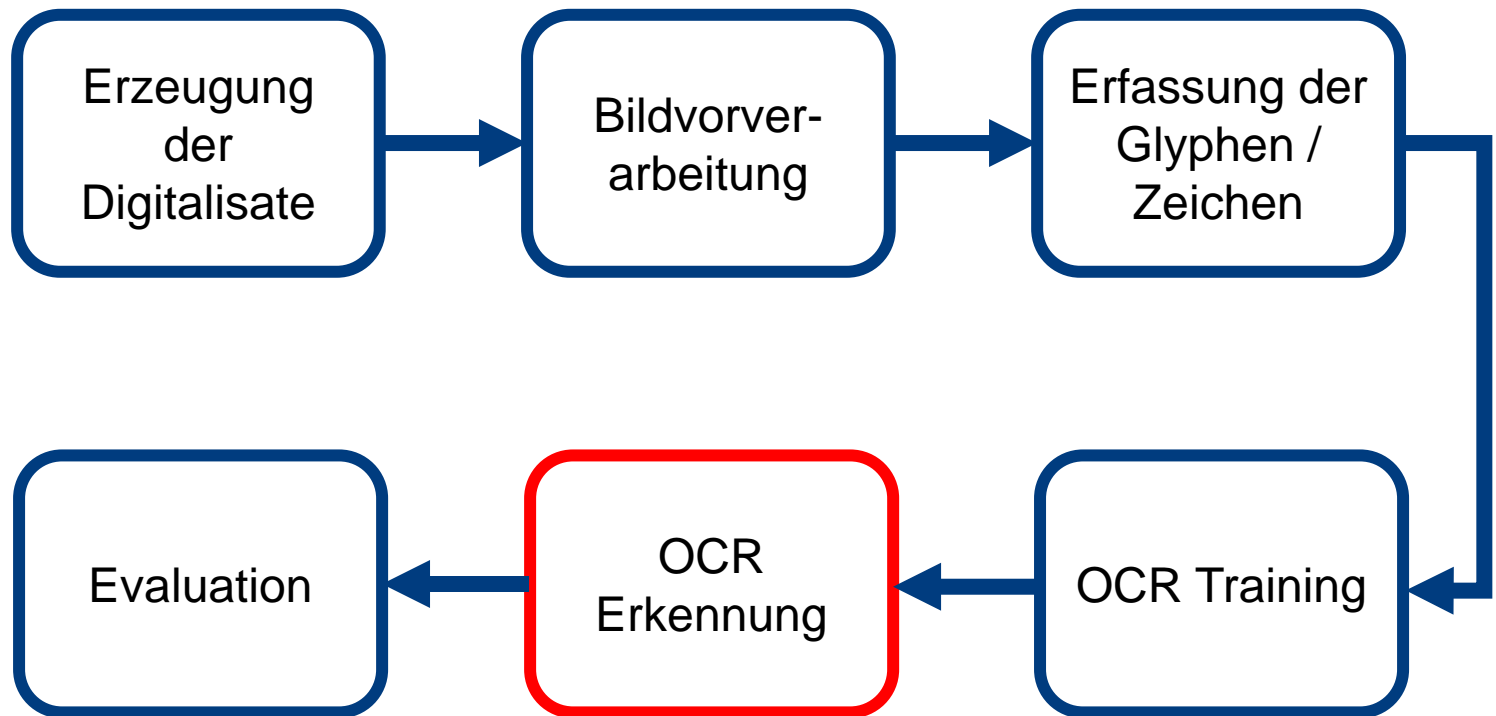
- PAGEXML
- *Ground Truth* in Textform

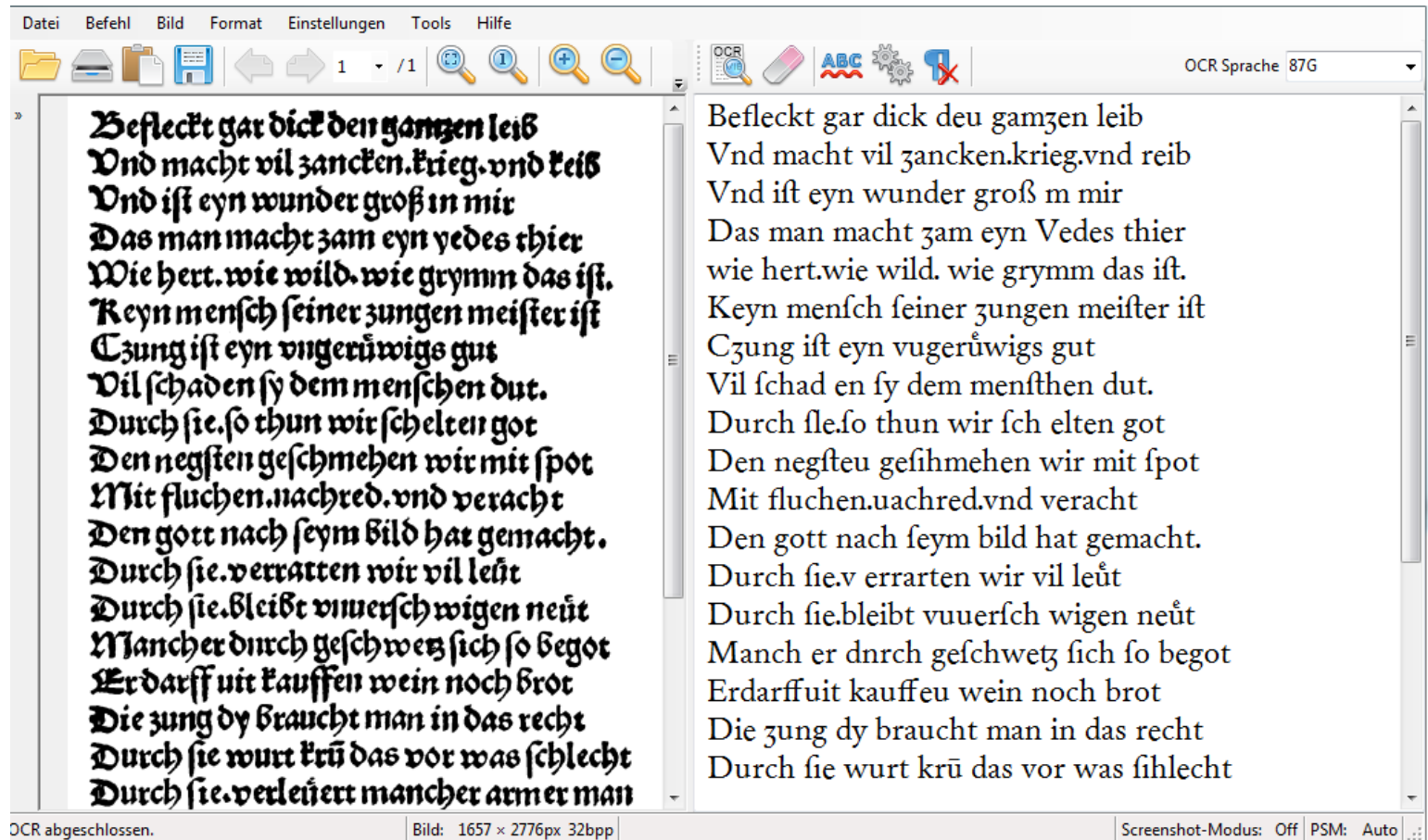


# Arbeitsschritte Franken+

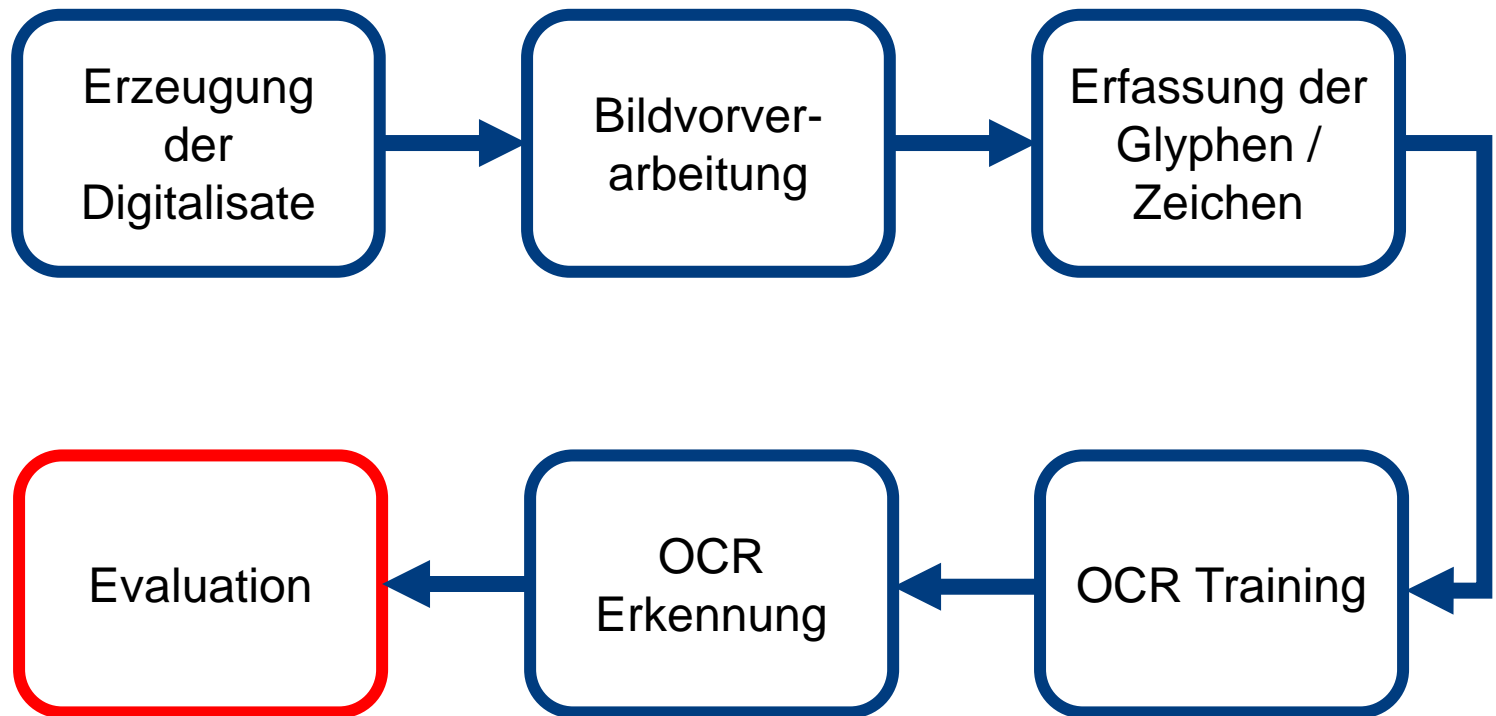
1. Font und Language erstellen
2. Glyphen (PAGEXML) importieren
3. Font bearbeiten
4. Synthetischen Text erstellen
5. Language trainieren

# OCR Workflow





# OCR Workflow





# TextEvaluation

PRImA Text Eval

About...

Ground Truth File  Select...

☐ UTF-8
 ☐ UTF-16
 ☐ UTF-16BE
 ☐ UTF-16LE
 ☐ US-ASCII
 ☐ ISO-8859-1

OCR Result File  Select...

☐ UTF-8
 ☐ UTF-16
 ☐ UTF-16BE
 ☐ UTF-16LE
 ☐ US-ASCII
 ☐ ISO-8859-1

Text Extraction Settings for XML Files (PAGE / ALTO)

☐ Additional line breaks
 ☐ Reading order only

Text Source 

region  
line  
word  
glyph

 Region Filter

Evaluation Method(s) 

BagOfWords  
CharacterAccuracy  
WordAccuracy

Settings

Stop Words 

EN  
DE

Convert to Lower Case 

DEFAULT  
CHINESE  
ENGLISH  
FRENCH  
GERMAN  
ITALIAN

Text Filter File  Select...

Evaluate Exit